

(12) INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(19) World Intellectual Property Organization
International Bureau



(43) International Publication Date
16 October 2003 (16.10.2003)

PCT

(10) International Publication Number
WO 03/085909 A1

(51) International Patent Classification⁷: **H04L 12/56,**
H04Q 11/04

Robert; 392 Clinton St. Apt. 2, Brooklyn, NY 11231 (US).

(21) International Application Number: PCT/US03/08218

(74) Agent: **WILLIAMS, Kirk**; The Law Office of Kirk D. Williams, (USPTO Customer No. 26327), 1234 S. Ogden St., Denver, CO 80210 (US).

(22) International Filing Date: 15 March 2003 (15.03.2003)

(25) Filing Language: English

(81) Designated States (*national*): AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DZ, EC, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ, NI, NO, NZ, OM, PH, PL, PT, RO, RU, SC, SD, SE, SG, SK, SL, TJ, TM, TN, TR, TT, TZ, UA, UG, UZ, VC, VN, YU, ZA, ZM, ZW.

(26) Publication Language: English

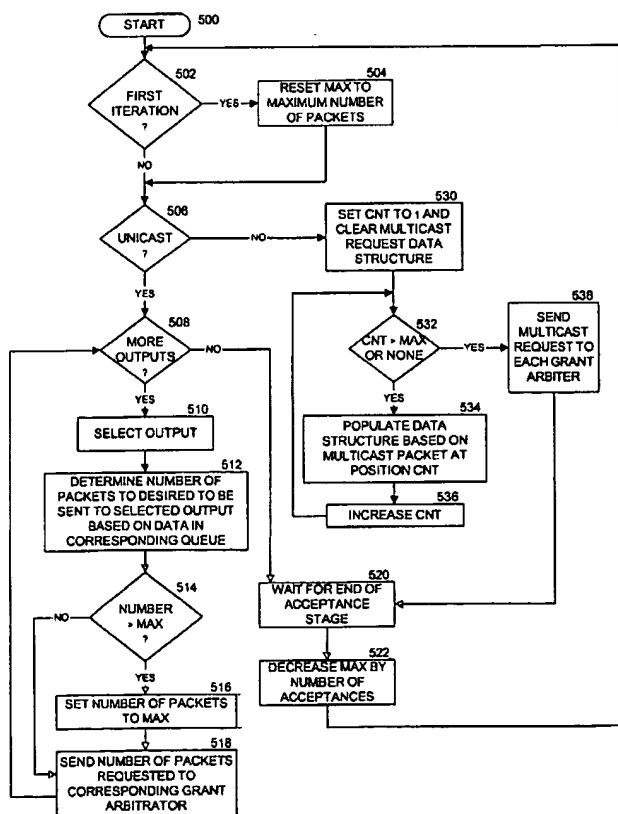
(30) Priority Data:
10/109,785 30 March 2002 (30.03.2002) US

(71) Applicant: **CISCO TECHNOLOGY, INC.** [US/US]; Patent Prosecution, c/o Robert Barr, 170 W. Tasman Dr., San Jose, CA 95134 (US).

(84) Designated States (*regional*): ARIPO patent (GH, GM, KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZM, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, BG, CH, CY, CZ, DE, DK, EE,

[Continued on next page]

(54) Title: PACKET SCHEDULING MECHANISM



(57) Abstract: Methods and apparatus are disclosed for scheduling packets, such as in systems having a non-blocking switching fabric and homogeneous or heterogeneous line card interfaces. In one implementation, multiple request generators, grant arbiters, and acceptance arbiters work in conjunction to determine this scheduling. A set of requests for sending packets from a particular input is generated. From a grant starting position, a first n requests in a predetermined sequence are identified, where n is less than or equal to the maximum number of connections that can be used in a single packet time to the particular output. The grant starting position is updated in response to the first n grants including a particular grant corresponding to a grant advancement position. In one embodiment, the set of grants generated based on the set of requests is similarly determined using an acceptance starting position and an acceptance advancement position.

WO 03/085909 A1



ES, FI, FR, GB, GR, HU, IE, IT, LU, MC, NL, PT, RO,
SE, SI, SK, TR), OAPI patent (BF, BJ, CF, CG, CI, CM,
GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG).

Published:

- with international search report
- before the expiration of the time limit for amending the claims and to be republished in the event of receipt of amendments

Declarations under Rule 4.17:

- as to applicant's entitlement to apply for and be granted a patent (Rule 4.17(ii)) for all designations
- as to the applicant's entitlement to claim the priority of the earlier application (Rule 4.17(iii)) for all designations

For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.

PACKET SCHEDULING MECHANISM

FIELD OF THE INVENTION

5 This invention especially relates to communications and computer systems; and more particularly, the invention relates to packet scheduling particularly applicable, but not limited to systems including a non-blocking switching fabric and homogeneous or heterogeneous line card interfaces.

10

BACKGROUND OF THE INVENTION

 The communications industry is rapidly changing to adjust to emerging technologies and ever increasing customer demand. This customer demand for new applications and increased performance of existing applications is driving communications network and system providers to employ networks and systems having greater speed and capacity (e.g., greater bandwidth). In trying to achieve these goals, a common approach taken by many communications providers is to use packet switching technology. Increasingly, public and private communications networks are being built and expanded using various packet technologies, such as Internet Protocol (IP).

20 SLIP is an iterative algorithm for scheduling the sending of packets across an $N \times N$ switch. In one implementation, the following three steps are performed:

1. Each unmatched input sends a request to every output for which it has a queued cell.
 2. If an unmatched output receives any requests, it chooses the one that appears next in a fixed, round-robin schedule starting from the highest selection priority element. The output notifies each input whether or not its request was granted. The pointer to the highest selection priority element of the round-robin schedule is incremented (modulo N) to one location beyond the granted input if and only if the grant is accepted in step 3 of the first iteration. The pointer is not incremented in subsequent iterations.
- 25
- 30

3. If an input receives a grant, it accepts the one that appears next in a fixed, round-robin schedule starting from the highest selection priority element. The pointer to the highest selection priority element of the round-robin schedule is incremented (modulo N) to one location beyond the accepted output.

5 I-SLIP is a scheduling algorithm including multiple iterations of the SLIP algorithm to determine the scheduling of packets for each round of sending packets (rather than just one SLIP iteration.)

Each output scheduler decides among the set of ordered, competing requests using a rotating selection priority. When a requesting input is granted and the input accepts that
10 grant, the input will have the lowest selection priority at that output in the next cell time. Also, whatever input has the highest selection priority at an output will continue to be granted during each successive time slot until it is serviced. This ensures that a connection will not be starved: the highest selection priority connection at an output will always be accepted by an input in no more than N cell times.

15 Moving the pointers not only prevents starvation, it tends to desynchronize the schedulers. Each of the outputs that matched in the previous time slot will have a different highest selection priority input. Thus, they will each grant to different inputs. Consider an example in which two inputs are both requesting the same two outputs. Initially, both outputs may grant to the same input; in that case only one connection will be made in the
20 first iteration.

The successful output will increment its pointer and in the next cell time, the outputs will no longer contend: one will have moved on to grant to another input and the other will grant to the same input as before. This leads to a better match in the first iteration of the next cell time. This is because the output schedulers have become desynchronized
25 (or "slipped") with respect to each other. This leads to high performance, even for a single iteration of SLIP.

Because of the round-robin movement of the pointers, the algorithm tends to provide a fair allocation of bandwidth among competing connections and to be burst-reducing. The burst-reduction is simplest to understand under high load when all input
30 queues are occupied: the algorithm will visit each competing connection in turn, so that

even if a burst of cells for the same output arrives at the input, the burst will be spread out in time if there is competing traffic.

However, the I-SLIP algorithm is designed to accommodate cross-bar switching fabrics wherein the input ports are independent and homogenous. Certain implementations of non-blocking switching fabrics have heterogeneous line cards of varying capacities. 5 Desired for these systems are schedulers that provide a reasonably fair bandwidth allocation across line cards of varying capacity, independently of the line card configuration. Even in systems wherein line cards of varying speeds are connected to a proportional increase in the number of input ports, the I-SLIP scheduling algorithm 10 typically does not provide a sufficiently fair bandwidth allocation. Needed are new methods and apparatus for scheduling packets across a non-blocking switching fabric and homogeneous or heterogeneous line card interfaces.

SUMMARY OF THE INVENTION

Methods and apparatus are disclosed for scheduling packets. In one embodiment, a set of requests corresponding to packets desired to be sent from a particular input across a packet switch to a plurality of outputs is generated. A grant starting position is maintained and a grant advancement position is determined. The first n requests in a predetermined sequence starting from the grant starting position are identified, where n is less than or equal to the maximum number of connections that can be used in a single packet time to the particular output. The grant starting position is updated in response to the first n grants if a particular grant corresponding to the grant advancement position is included.

In one embodiment, a set of requests corresponding to packets desired to be sent from a particular input across a packet switch to a plurality of outputs is generated, and a set of grants generated in response to the set of requests are identified. An acceptance starting position is maintained, and an acceptance advancement position is determined. The first m grants in a predetermined sequence starting from the acceptance starting position are identified, where m is less than or equal to the maximum number of connections that can be used in a single packet time from the particular input. The acceptance starting position is updated in response to the first m grants if a particular grant corresponding to the acceptance advancement position is included.

20

25

BRIEF DESCRIPTION OF THE DRAWINGS

The appended claims set forth the features of the invention with particularity. The invention, together with its advantages, may be best understood from the following detailed description taken in conjunction with the accompanying drawings of which:

5 FIGs. 1A-E and 2 are block diagrams of embodiments scheduling packets in a system having a non-blocking switching fabric;

 FIG. 3A is a flow diagram of a process used in one embodiment for scheduling unicast and multicast packets in three iteration scheduling cycles;

10 FIG. 3B is a flow diagram of a process used in one embodiment for scheduling unicast and/or multicast packets in one or more iterations;

 FIGs. 4A and 4C are flow diagrams of processes used in one embodiment for communicating unicast and multicast packet indications to a scheduler;

 FIG. 4B is a block diagram of a message format used in one embodiment for communicating unicast and multicast packet indications to a scheduler;

15 FIG. 5 is a flow diagram of a process used in one embodiment for generating requests;

 FIG. 6A is a flow diagram of a process used in one embodiment in performing grant processing;

20 FIGs. 6B-C are block diagrams of data structures used in one embodiment in performing grant processing;

 FIG. 7A is a flow diagram of a process used in one embodiment for performing acceptance processing;

 FIGs. 7B illustrates block diagrams of data structures used in one embodiment for performing acceptance processing;

25 FIG. 8 is a flow diagram of a process used in one embodiment for multicast pointer processing; and

 FIG. 9 is a block diagram used in one embodiment for configuring the switch and initiating the sending of packets across the switch.

DETAILED DESCRIPTION

Methods and apparatus are disclosed for scheduling packets in systems, such as, but not limited to systems having a non-blocking switching fabric and homogeneous or heterogeneous line card interfaces. Embodiments described herein include various elements and limitations, with no one element or limitation contemplated as being a critical element or limitation. Each of the claims individually recites an aspect of the invention in its entirety. Moreover, some embodiments described may include, but are not limited to, *inter alia*, systems, networks, integrated circuit chips, embedded processors, ASICs, methods, and computer-readable medium containing instructions. The embodiments described hereinafter embody various aspects and configurations within the scope and spirit of the invention, with the figures illustrating exemplary and non-limiting configurations.

As used herein, the term "packet" refers to packets of all types or any other units of information or data, including, but not limited to, fixed length cells and variable length packets, each of which may or may not be divisible into smaller packets or cells. The term "packet" as used herein also refers to both the packet itself or a packet indication, such as, but not limited to all or part of a packet or packet header, a data structure value, pointer or index, or any other part or identification of a packet. Moreover, these packets may contain one or more types of information, including, but not limited to, voice, data, video, and audio information. The term "item" is used herein to refer to a packet or any other unit or piece of information or data.

The term "system" is used generically herein to describe any number of components, elements, sub-systems, devices, packet switch elements, packet switches, routers, networks, computer and/or communication devices or mechanisms, or combinations of components thereof. The term "computer" is used generically herein to describe any number of computers, including, but not limited to personal computers, embedded processors and systems, control logic, ASICs, chips, workstations, mainframes, etc. The term "device" is used generically herein to describe any type of mechanism, including a computer or system or component thereof. The terms "task" and "process" are used generically herein to describe any type of running program, including, but not limited to a computer process, task, thread, executing application, operating system, user process,

device driver, native code, machine or other language, etc., and can be interactive and/or non-interactive, executing locally and/or remotely, executing in foreground and/or background, executing in the user and/or operating system address spaces, a routine of a library and/or standalone application, and is not limited to any particular memory partitioning technique. The steps, connections, and processing of signals and information illustrated in the figures, including, but not limited to any block and flow diagrams and message sequence charts, may be performed in the same or in a different serial or parallel ordering and/or by different components and/or processes, threads, etc., and/or over different connections and be combined with other functions in other embodiments in keeping within the scope and spirit of the invention.

Moreover, the terms "network" and "communications mechanism" are used generically herein to describe one or more networks, communications mediums or communications systems, including, but not limited to the Internet, private or public telephone, cellular, wireless, satellite, cable, local area, metropolitan area and/or wide area networks, a cable, electrical connection, bus, etc., and internal communications mechanisms such as message passing, interprocess communications, shared memory, etc.

The term "storage mechanism" includes any type of memory, storage device or other mechanism for maintaining instructions or data in any format. "Computer-readable medium" is an extensible term including any memory, storage device, storage mechanism, and other storage and signaling mechanisms including interfaces and devices such as network interface cards and buffers therein, as well as any communications devices and signals received and transmitted, and other current and evolving technologies that a computerized system can interpret, receive, and/or transmit. The term "memory" includes any random access memory (RAM), read only memory (ROM), flash memory, integrated circuits, and/or other memory components or elements. The term "storage device" includes any solid state storage media, disk drives, diskettes, networked services, tape drives, and other storage devices. Memories and storage devices may store computer-executable instructions to be executed a processor and/or control logic, and data which is manipulated a processor and/or control logic. The term "data structure" is an extensible term referring to any data element, variable, data structure, data base, and/or one or more or an organizational schemes that can be applied to data to facilitate interpreting the data or

performing operations on it, such as, but not limited to memory locations or devices, sets, queues, trees, heaps, lists, linked lists, arrays, tables, pointers, etc. A data structure is typically maintained in a storage mechanism.

The terms "first," "second," etc. are typically used herein to denote different units
5 (e.g., a first element, a second element). The use of these terms herein does not necessarily connote an ordering such as one unit or event occurring or coming before the another, but rather provides a mechanism to distinguish between particular units. Moreover, the phrases "based on x" and "in response to x" are used to indicate a minimum set of items x from
10 which something is derived or caused, wherein "x" is extensible and does not necessarily describe a complete list of items on which the is performed, etc. Additionally, the phrase "coupled to" is used to indicate some level of direct or indirect connection between two elements or devices, with the coupling device or devices modify or not modifying the coupled signal or communicated information. The term "subset" is used to indicate a group of all, less than all, or none of the elements of a set. Moreover, the term "or" is used herein
15 to identify an alternative selection of one or more, including all, of the conjunctive items.

Methods and apparatus are disclosed for scheduling packets in systems, such as, but not limited to systems having a non-blocking switching fabric and homogeneous or heterogeneous line card interfaces. In one embodiment, multiple request generators, grant arbiters, and acceptance arbiters work in conjunction to determine this scheduling. A set of
20 requests for sending packets from a particular input is generated. From a grant starting position, a first n requests in a predetermined sequence are identified, where n is less than or equal to the maximum number of connections that can be used in a single packet time to the particular output. The grant starting position is updated in response to the first n grants including a particular grant corresponding to a grant advancement position. In one
25 embodiment, the set of grants generated based on the set of requests is similarly determined using an acceptance starting position and an acceptance advancement position.

In one embodiment, a "packet time" is a time interval for a given switch configuration during which one or more packets can be sent from one or more inputs to one or more outputs. In one embodiment, the packet time corresponds to the scheduling time
30 interval required or allocated to perform the scheduling of packets, and thus, packets can be

sent while the packet scheduling and corresponding switch configuration are being determined for the next packet time.

FIG. 1 illustrates one embodiment of a system 100 including a non-blocking switch (or switch fabric) 102, a control with scheduler and memory 101, and multiple line cards 103-106. Line card 103 is denoted as being of "type A" with A1 ingress links or ports 104 and A2 egress links or ports 105. Line card 106 is denoted as being of "type B" with N1 ingress links or ports 107 and N2 egress links or ports 108. This labeling emphasizes that interfaces and line cards with varying rates and numbers of ports or connections to a non-blocking switch 102 are supported.

FIG. 1B illustrates one embodiment of a line card 110. Signals including packets or other data formats are received and transmitted by line interface 111. Shown are unicast and multicast queues 113, wherein incoming packets to be scheduled are placed in one embodiment. Control with request generators, grant arbiters, and acceptance arbiters 112 determines and schedules packets as described hereinafter, with packets being sent from unicast and multicast queues 113 at their respective scheduled times via switch interface 114. Additionally, scheduling requests, grants, and acceptances are communicated among other request generators, grant arbiters, and acceptance arbiters via switch interface 114.

FIG. 1C illustrates one embodiment wherein the request generators, grant arbiters, and acceptance arbiters are centrally located in control with request generators, grant arbiters and acceptance arbiters 122. Line cards with unicast and multicast queues and packet indication generators 121 send packet traffic indications 123 to control with request generators, grant arbiters and acceptance arbiters 122. Returned are acceptance/schedule indications 124 of packets to line cards 121, which initiate the sending of the accepted packets at the scheduled time. Additionally, control with request generators, grant arbiters and acceptance arbiters 122 sends configuration information 125 to switch 120, so the switching fabric can be configured to communicate the accepted packets between the switch input and output ports and connected line cards 121.

FIG. 1D illustrates one embodiment of a line card 130. Signals including packets or other data formats are received and transmitted by line interface 131. Shown are N unicast queues 133-134 and one multicast queue 135, wherein incoming packets to be scheduled are placed. Typically, N corresponds to the number of output line cards or the number of

switch output ports to which the line card can send packets. In one embodiment, additional queues are used, such as, but not limited to multiple multicast queues and queues for buffering packets having various priority levels. Control with request module and memory 132 sends packet indications and receives acceptance and scheduling indications via switch interface 136.

FIG. 1E illustrates a system 150 including a N request generators 154, grant arbiters 155, and acceptance arbiters 156. Packet indications are received from various line cards via switch interface 151 and stored in the corresponding queue of the N unicast queues 152 and N multicast request queues 153. The N request generators 154, based on the packet indications in queues 152 and 153, generate unicast and multicast packet requests (typically in separate iterations) and communicate to the grant arbiters corresponding to the destination of the packets of the N grant arbiters 155. The N grant arbiters 155 in turn generate and communicate their grants to the acceptance arbiters corresponding to the source of the granted packets of the N acceptance arbiters 156. The acceptances are then, or after multiple iterations, communicated to switch interface 151 for relaying to the appropriate line cards and switch configuration control. In one embodiment, a multicast control 157 is used maintain a common multicast position used by grant arbiters 155 in selecting which multicast requests to grant.

FIG. 2 illustrates one embodiment of a system 200, which may include, but is not limited to one or more request generators, grant arbiters and/or acceptance arbiters for scheduling packets according to the invention. In one embodiment, system 200 includes a processor 201, memory 202, storage devices 203, and switch/control interface 204, which are typically coupled via one or more communications mechanisms 209 (shown as a bus for illustrative purposes.) Various embodiments of system 200 may include more or less elements. The operation of system 200 is typically controlled by processor 201 using memory 202 and storage devices 203 to perform one or more scheduling tasks or processes. Memory 202 is one type of computer-readable medium, and typically comprises random access memory (RAM), read only memory (ROM), flash memory, integrated circuits, and/or other memory components. Memory 202 typically stores computer-executable instructions to be executed by processor 201 and/or data which is manipulated by processor 201 for implementing functionality in accordance with the

invention. Storage devices 203 are another type of computer-readable medium, and typically comprise solid state storage media, disk drives, diskettes, networked services, tape drives, and other storage devices. Storage devices 203 typically store computer-executable instructions to be executed by processor 201 and/or data which is manipulated by processor 201 for implementing functionality in accordance with the invention.

FIG. 3A illustrates a process used in one embodiment for scheduling packets using three scheduling iterations. Processing begins in process block 300, and proceeds to process block 302, wherein a first unicast scheduling iteration is performed. Next, in process block 304, a second unicast scheduling iteration is performed. In process block 306, a multicast scheduling iteration is performed. Next, in process block 308, the switch (and its switching fabric) are configured according to the scheduled packets, and in process block 310, the packets are sent. For the next scheduling round, processing proceeds to process block 312, wherein a multicast scheduling iteration is performed. Next, in process block 314, a first unicast scheduling iteration is performed. In process block 316, a second unicast scheduling iteration is performed. Next, in process block 318, the switch (and its switching fabric) are configured according to the scheduled packets, and in process block 319, the packets are sent. Processing returns to process block 302 to perform more scheduling of packets.

FIG. 3B illustrates a process used in one embodiment for scheduling packets using one or more scheduling iterations, including unicast and/or multicast iterations in any desired order. Processing begins with process block 320. As determined in process block 322, if a unicast iteration is next, then in process block 324, the unicast scheduling iteration is performed; otherwise, a multicast scheduling iteration is performed in process block 326. As determined in process block 328, if there are more scheduling iterations to be performed for this scheduling cycle, then processing returns to process block 322 to perform the next scheduling iteration. Otherwise, the switch is configured in process block 330, packets are sent in process block 340, and processing then returns to process block 322.

FIG. 4A illustrates a process used in one embodiment to generate packet indication messages. Processing begins with process block 400, and proceeds to process block 402, wherein a packet indication data structure is cleared. As determined in process block 404,

if there are more unicast packets to be sent, then a first or next position in the unicast queues is selected in process block 406. In process block 408, a bitmap or other representation of the destination or destinations of the packets at the selected position in the destination queues is added to the data structure, and processing returns to process block 404. In one embodiment for unicast and/or multicast packets, if a particular destination is disabled, out of service, or currently unreachable based on backpressure or other flow control information, indications for this destination are not added to the data structure in process blocks 408 or 414.

Otherwise, as determined in process block 410, if there are more multicast packets to be sent, then a first or next position in the multicast queue is selected in process block 412. In process block 414, a bitmap or other representation of the destinations of the multicast packet at the selected position in the multicast queue is added to the data structure, and processing returns to process block 410.

Otherwise, the data structure is sent to the scheduler in process block 430. In process block 432, indications are received from the scheduler of which packets to send and the multicast queues are updated if less than all destinations of a particular packet are allowed. The sending of these packets is initiated in process block 434. Processing returns to process block 402.

FIG. 4B illustrates a block diagram of a data structure/message format 450 used in one embodiment. Data structure 450 typically has multiple entries, each with an identification field 451 to indicate whether the entry corresponds to unicast or multicast packet indications, and a bitmap field 452 to indicate the destinations of the packets.

FIG. 4C illustrates a process used in one embodiment by a centralized scheduling system to collect the packet indications for the various sending line cards. Processing begins with process block 470, and proceeds to process block 472, wherein a message is received. In process block 474, one or more packet indication queues or other data structures are updated, and processing returns to process block 472.

FIG. 5 illustrates a process used in one embodiment by each of the request generators, typically one for each line card associated with the non-blocking packet switch. Processing begins with process block 500. As determined in process block 502, if this is a first iteration, then in process block 504, the value of MAX is set to the maximum number

of packets that can be sent by the line card in a packet time, which typically corresponds to the number of switch input ports to which the line card connects. Each request generator will typically have outstanding a cumulative number of requests that it can service in a scheduling cycle.

5 As determined in process block 506, if this is a unicast iteration, then processing proceeds to process block 508 to indicate a set of requests to each of the grant arbiters. While there are more outputs as determined in process block 508, an output is selected in process block 510, and the number of desired packets to be sent to the particular output (up to the maximum number of packets the destination can actually receive in a packet time) is
10 determined in process block 512. If this number is greater than the value of MAX as determined in process block 514, then this number is set to MAX in process block 516. In process block 518, the requests are signaled to the corresponding grant arbiter. After all outputs have been processed, then in process block 520, the request arbiter waits for the end of the acceptance stage of the current unicast iteration. Then, in process block 522,
15 MAX is decreased by the number of acceptances corresponding to the previously sent requests from this request arbiter in this iteration, and processing returns to process block 502.

 If, as determined in process block 506, that this is a multicast iteration, then processing proceeds to process block 530 to set CNT to one and to clear the multicast
20 request data structure. While CNT is not greater than MAX and there are multicast requests to process as determined in process block 532, processing blocks 534 and 536 are performed. In process block 534, a data structure is populated based on the destinations of the multicast packet at position CNT in the multicast queue, and CNT is increased by one in process block 536. When done, processing proceeds to process block 538 to send a
25 multicast request to each grant arbiter (of course, it could be a request of no multicast packets) or at least those grant arbiters with a pending multicast request from this request generator. Processing then proceeds to process block 520.

 FIG. 6A illustrates a flow diagram of a process used by a grant arbiter in one embodiment. Processing begins with process block 600, and proceeds to process block
30 602, wherein a grant starting position is initialized. Next, in process block 604, the requests are received from the request generators, with these requests used to populate a data

structure. In one embodiment, data structure 650 illustrated in FIG. 6B is used, with data structure 650 including a bitmap unary representation of the number of requests received for each slot (e.g., from each request generator).

5 In one embodiment, these bitmap representations are right-aligned as illustrated in data structure 660. In one embodiment, these bitmap representations are left-aligned, while in one embodiment, these bitmap representations are maintained in a manner representative of the physical ports of the line card or slot. The alignment of the requesting bits within such a bitmap typically impacts packet scheduling by affecting the updating of the grant starting position. When the bitmap is right-aligned, the starting position for selecting bits
10 (e.g., bits corresponding to grants or acceptances) is more likely to advance to bits corresponding to a next line card or slot. However, this rate of advancement is still throttled by, *inter alia*, the traffic rate of the line card and switch throughput as indicated by the generation rate of requests, grants, and acceptances, as well as the line cards and ports corresponding to the particular requests, grants, and acceptances.

15 Returning to the processing of FIG. 6A, as determined in process block 606, if this is a first iteration of the current scheduling round, then in process block 608, MAX is set to the maximum number of packets which can be received in one packet time by the line card corresponding to this grant arbiter. Next, as determined in process block 610, if this is a unicast iteration, then in process block 612, the grant advancement position (GAP) is
20 determined. If a grant corresponding to the grant advancement position is accepted during the first iteration (or in any iteration in one embodiment), then the grant starting position will be modified so grants will be generated starting from a different position in a next scheduling round.

In one embodiment, the grant advancement position is the first position in the
25 request data structure indicating a request after the grant starting position. Referring back to FIG. 6C, data structure 660 illustrates two right-aligned bitmaps. If the grant starting position is at position 661, then the grant advancement position is at position 662. If the grant starting position is at position 662, then the grant advancement position is at position 663. If the grant starting position is at position 663, then the grant advancement position is
30 at position 664.

Returning to the processing of FIG. 6A and process block 614, if the iteration is not a unicast iteration, then in process block 616, up to MAX multicast requests are generated beginning at the multicast pointer position (common among all grant arbiters in one embodiment), and these grants are sent to the corresponding acceptance arbiters.

5 Otherwise, in process block 618, up to MAX unicast grants are generated beginning at the grant starting position. Next, in process block 620, these generated grants, along with an indication of whether a grant at the grant advancement position is included, are sent to the corresponding acceptance arbiters. Next, in process block 622, indications of the accepted grants are received, and MAX is decreased by the number of accepted grants
10 generated by this grant arbiter. If, as determined in process block 624, this is a first iteration of the current scheduling cycle, then as determined in process block 626, if the packet at the grant advancement position was accepted, then the advance flag is set in process block 628. As determined in process block 630, if this is a last iteration of the current scheduling cycle, then as determined in process block 632, if the advance flag is set, then in process
15 block 634, the grant starting position is advanced to the next position after the grant advancement position. Processing then returns to process block 604.

FIG. 7A illustrates a flow diagram of a process used by an acceptance arbiter in one embodiment. Processing begins with process block 700, and proceeds to process block 702, wherein an acceptance starting position is initialized. Next, in process block 704, the
20 grants and grant advancement position indicators are received from the grant arbiters, with this data being used to populate one or more data structures. In one embodiment, GAP data structure 740 illustrated in FIG. 7B is used to maintain the grant acceptance indications for each of the grant arbiters (corresponding to line card slots in one embodiment), and grant data structure 750 including a bitmap unary representation of the number of grants received
25 for each slot (e.g., from each request generator). These bitmaps may or may not be right-aligned.

Returning to the processing of FIG. 7A and process block 706, if this is a unicast iteration and a first iteration of the scheduling cycle, then in process block 708, the acceptance advancement position is typically determined in the same manner as that for the
30 grant advancement position as described herein.

Next, as determined in process block 710, if this is a multicast iteration, then in process block 712, all grants are accepted (as a sending line card does not send more multicast requests than it can service), acceptance indications are transmitted, and processing returns to process block 704.

- 5 Otherwise, in process block 714, up to MAX unicast grants are accepted beginning with the grant at the acceptance advancement position, then grants from the grant starting position. Next, in process block 716, the corresponding grant arbiters are notified of their accepted grants and whether their GAP grant was accepted. Next, in process block 718, MAX is decreased by the number of accepted grants generated by this acceptance arbiter.
- 10 If, as determined in process block 720, this is a first iteration of the current scheduling cycle, then as determined in process block 722, if the grant at the acceptance advancement position was accepted, then the advance flag is set in process block 724. As determined in process block 726, if this is a last iteration of the current scheduling cycle, then as determined in process block 728, if the advance flag is set, then in process block 730, the
- 15 acceptance starting position is advanced to the next position after the acceptance advancement position. Processing then returns to process block 704.

FIG. 8 illustrates a process used in one embodiment by a multicast control to update the multicast pointer. Processing begins at process block 800, and proceeds to process block 802, wherein the multicast starting position is initialized. Next, in process block 804,

20 multicast request messages are received from the various request generators. In process block 806, the multicast advancement position is set to the next position having a multicast request at or after the multicast starting position. In process block 808, multicast acceptance indications are received. As determined in process block 810, if all the requests for the multicast packet at the head of the queue corresponding to the multicast starting

25 position were accepted (e.g., the first multicast packet to be sent from the input corresponding to the MAP position was fully accepted), then in process block 812, the multicast starting position is set to the next position after the multicast advancement position. Processing returns to process block 804.

FIG. 9 illustrates a process used in one embodiment for configuring a switch (e.g.,

30 non-blocking switch fabric) and sending of the accepted packets. Processing begins with process block 900, and proceeds to process block 902, wherein indications of the accepted

connection are received. In process block 904, the switch is configured at the appropriate time to connect the appropriate input and output ports of the switch corresponding to the accepted requests. Then, in process block 906, sending of the packets are initiated and sent. Processing returns to process block 902.

- 5 In view of the many possible embodiments to which the principles of our invention may be applied, it will be appreciated that the embodiments and aspects thereof described herein with respect to the drawings/figures are only illustrative and should not be taken as limiting the scope of the invention. For example and as would be apparent to one skilled in the art, many of the process block operations can be re-ordered to be performed before,
- 10 after, or substantially concurrent with other operations. Also, many different forms of data structures could be used in various embodiments. The invention as described herein contemplates all such embodiments as may come within the scope of the following claims and equivalents thereof.

CLAIMS

What is claimed is:

1. A method, comprising:
 - identifying a set of requests corresponding to packets desired to be sent from a
 - 5 plurality of inputs across a packet switch to a particular output;
 - maintaining a grant starting position;
 - determining a grant advancement position;
 - identifying a first n requests in a predetermined sequence starting from the grant
 - starting position, where n is less than or equal to the maximum number of packets that can
 - 10 be sent in a single packet time to the particular output; and
 - updating the grant starting position in response to the first n grants including a
 - particular grant corresponding to the grant advancement position.
2. The method of claim 1, wherein said updating the grant starting position is
- performed only if the first n grants includes the particular grant corresponding to the grant
- 15 advancement position.
3. The method of claim 1, wherein said updating the acceptance starting position
- includes setting a value of the grant starting position to one position beyond the grant
- advancement position.
4. The method of claim 1, further comprising:
- 20 maintaining a bitmap data structure indicating the set of requests; and
- wherein the grant starting position and the grant advancement position identify one
- or more positions within the bitmap data structure.
5. The method of claim 4, wherein the bitmap data structure includes a bitmap
- representation corresponding to each of the plurality of inputs.
- 25 6. The method of claim 5, wherein the bitmap representations are right-aligned.
7. The method of claim 1, wherein the set of requests includes at most n requests
- from each input of the plurality of inputs.

8. The method of claim 1, wherein the set of request includes from each particular input of the plurality of inputs at most the maximum number of packets that said each particular input can send in the single packet time.

9. The method of claim 8, wherein the set of requests includes at most n requests
5 from each input of the plurality of inputs.

10. A computer-readable medium containing computer-executable instructions for performing the method of claim 1.

11. A method, comprising:
generating a set of requests corresponding to packets desired to be sent from a
10 particular input across a packet switch to a plurality of outputs;
identifying a set of grants generated in response to the set of requests;
maintaining an acceptance starting position;
determining an acceptance advancement position;
identifying a first m grants in a predetermined sequence starting from the
15 acceptance starting position, where m is less than or equal to the maximum number of
connections that can be used in a single packet time from the particular input; and
updating the acceptance starting position in response to the first m grants including
a particular acceptance corresponding to the acceptance advancement position.

12. The method of claim 11, wherein said updating the acceptance starting position
20 is performed only if the first m grants includes the particular acceptance corresponding to
the acceptance advancement position.

13. The method of claim 11, wherein said updating the acceptance starting position
includes setting a value of the acceptance starting position to one position beyond the
acceptance advancement position.

25 14. The method of claim 11, further comprising:
maintaining a bitmap data structure indicating the set of grants; and
wherein the acceptance starting position and the acceptance advancement position
identify one or more positions within the bitmap data structure.

15. The method of claim 14, wherein the bitmap data structure includes a bitmap representation corresponding to each of the plurality of outputs.

16. The method of claim 15, wherein the bitmap representations are right-aligned.

17. The method of claim 11, wherein the set of requests includes at most m requests
5 from the particular input.

18. The method of claim 11, wherein the set of request includes for each particular output of the plurality of outputs at most the maximum number of packets that can be sent to said each particular output in the single packet time.

19. The method of claim 18, wherein the set of requests includes at most m requests
10 from the particular input.

20. A computer-readable medium containing computer-executable instructions for performing the method of claim 11.

21. An apparatus, comprising:

a plurality of request generators;

a plurality of grant arbiters coupled to the plurality of request generators;

a plurality of acceptance arbiters coupled to the plurality of grant arbiters;

5 wherein each of the plurality of request generators is configured for generating requests for its associated input of a plurality of inputs of a switch;

 wherein each of the plurality of grant arbiters is configured for generating grants based on one or more received requests, said grants corresponding to authorization to send to its associated output of a plurality of outputs of the switch, where said generating grants
10 includes maintaining a grant starting position, determining a grant advancement position, identifying a first n requests in a predetermined sequence starting from the grant starting position, where n is less than or equal to the maximum number of packets that can be sent in a single packet time to said associated output; and updating the grant starting position in response to the first n grants including a particular grant corresponding to the grant
15 advancement position; and

 wherein each of the plurality of acceptance arbiters is configured for generating acceptances based on one or more received grants, said acceptances corresponding to its associated input of the plurality of inputs.

22. The apparatus of claim 21, further comprising a plurality of packet indication
20 generators, each of the plurality of packet indication generators associated with a different input of the plurality of inputs and a different one of the plurality of request generators;

 wherein each of the plurality of packet indication generators is configured to send an indication of queued packets and their respective destinations to its associated one of the plurality of request generators.

23. The apparatus of claim 21, wherein packets are sent across the switch from the
25 plurality of inputs to a plurality of outputs of the switch according to said acceptances.

24. The apparatus of 23, wherein said acceptances include a first unicast set of acceptances, a second unicast set of acceptances, and a multicast set of acceptances.

25. The apparatus of claim 21, wherein said generating acceptances includes maintaining an acceptance starting position, determining an acceptance advancement position, identifying a first m grants in a predetermined sequence starting from the acceptance starting position, where m is less than or equal to the maximum number of connections that can be used in a single packet time from said associated input, and updating the acceptance starting position in response to the first m grants including a particular acceptance corresponding to the acceptance advancement position.

26. The apparatus of claim 21, further comprising a multicast control configured to maintain a multicast pointer.

27. The apparatus of claim 26, wherein the multicast control is further configured to determine a multicast advancement position, and to update the multicast pointer in response to a particular request corresponding to the multicast advancement position being granted, a particular grant corresponding to the multicast advancement position being accepted, or a multicast packet corresponding to the multicast advancement position being sent.

28. An apparatus, comprising:
means for identifying a set of requests corresponding to packets desired to be sent from a plurality of inputs across a packet switch to a particular output;
means for maintaining a grant starting position;
means for determining a grant advancement position;
means for identifying a first n requests in a predetermined sequence starting from the grant starting position, where n is less than or equal to the maximum number of packets that can be sent in a single packet time to the particular output; and
means for updating the grant starting position in response to the first n grants including a particular grant corresponding to the grant advancement position.

29. The apparatus of claim 28, wherein said means for updating the acceptance starting position includes means for setting a value of the grant starting position to one position beyond the grant advancement position.

30. The apparatus of claim 28, further comprising:

means for maintaining a bitmap data structure indicating the set of requests; and
wherein the grant starting position and the grant advancement position identify one or more positions within the bitmap data structure.

5 31. An apparatus, comprising:

means for generating a set of requests corresponding to packets desired to be sent
from a particular input across a packet switch to a plurality of outputs;

means for identifying a set of grants generated in response to the set of requests;

means for maintaining an acceptance starting position;

10 means for determining an acceptance advancement position;

means for identifying a first m grants in a predetermined sequence starting from the
acceptance starting position, where m is less than or equal to the maximum number of
packets that can be sent in a single packet time from the particular input; and

means for updating the acceptance starting position in response to the first m grants

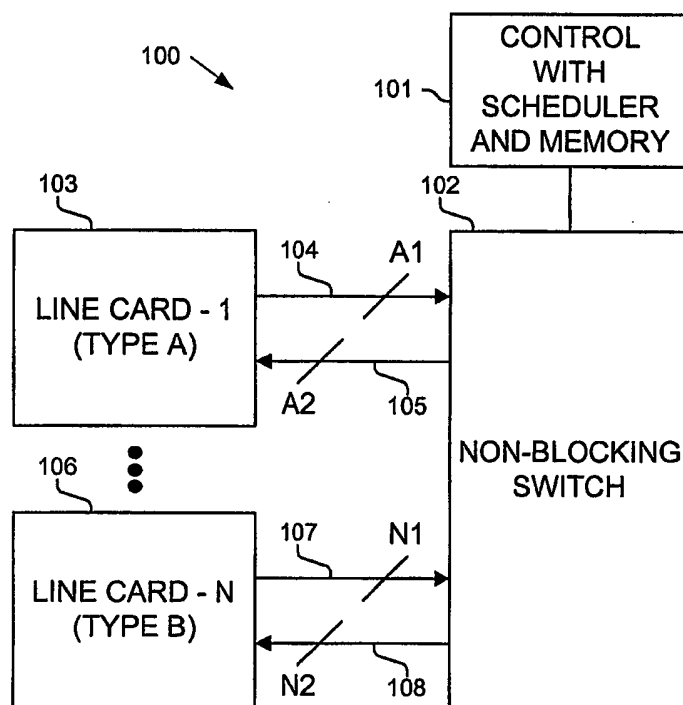
15 including a particular acceptance corresponding to the acceptance advancement position.

32. The apparatus of claim 31, wherein said means for updating the acceptance
starting position includes means for setting a value of the acceptance starting position to
one position beyond the acceptance advancement position.

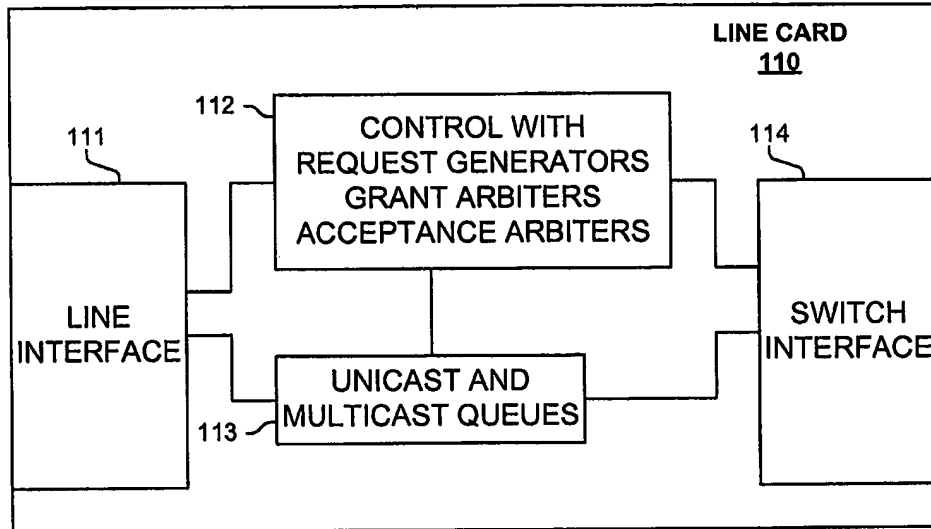
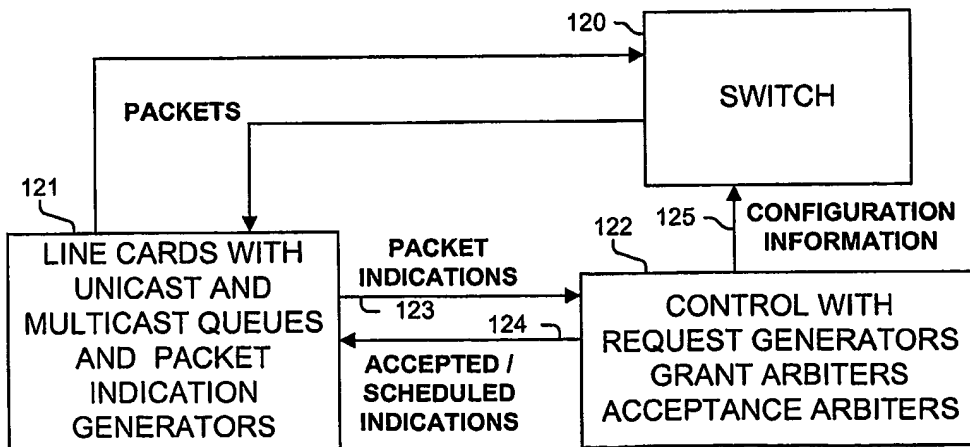
33. The apparatus of claim 31, further comprising:

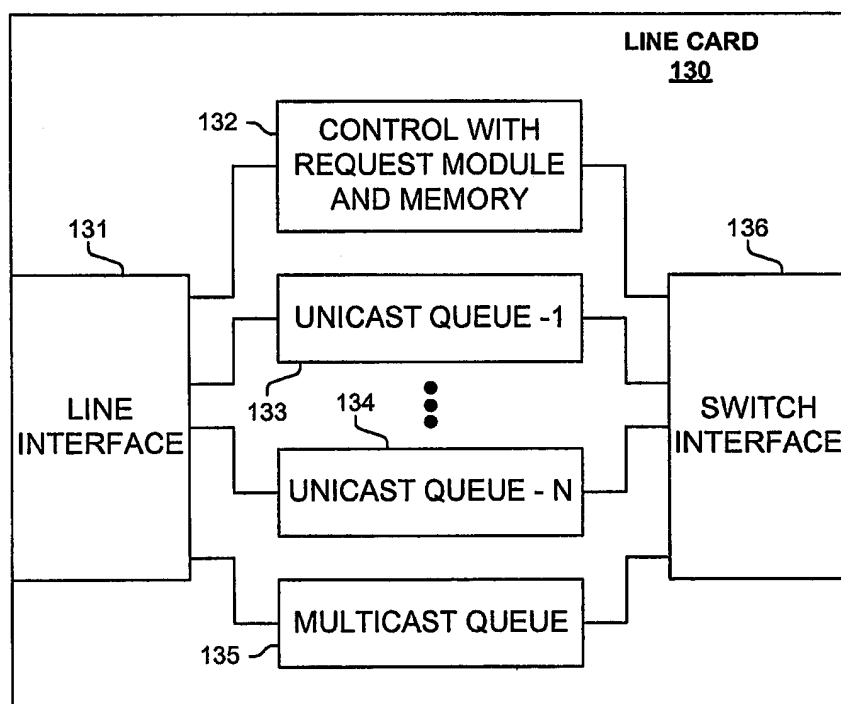
20 means for maintaining a bitmap data structure indicating the set of grants; and

wherein the acceptance starting position and the acceptance advancement position
identify one or more positions within the bitmap data structure.

**FIGURE 1A**

2 / 16

**FIGURE 1B****FIGURE 1C**

**FIGURE 1D**

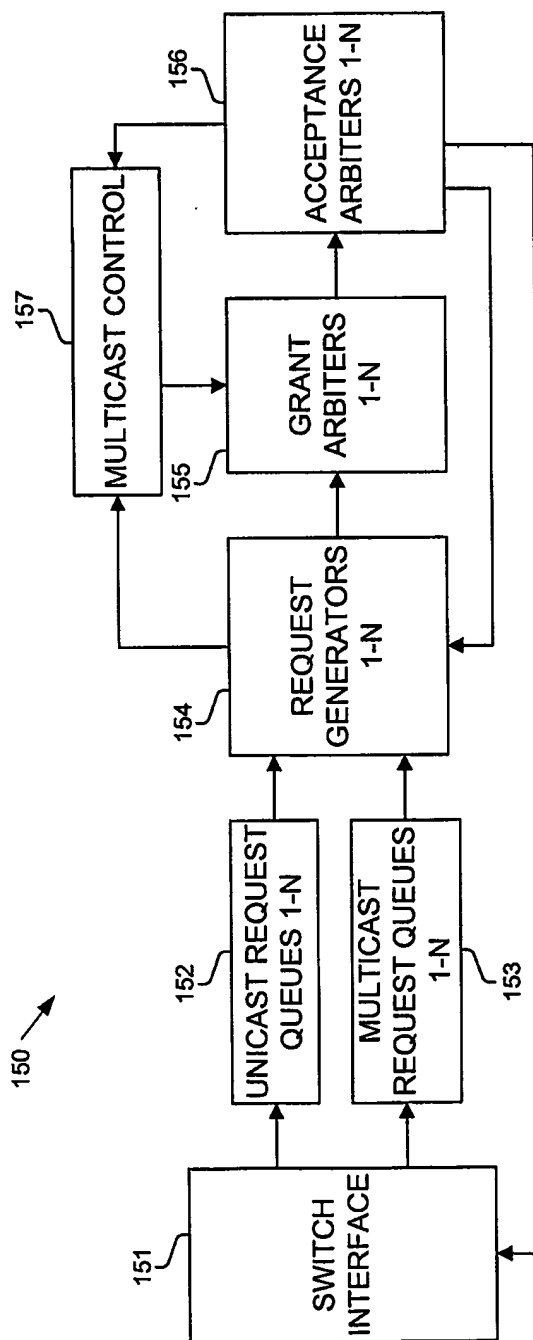


FIGURE 1E

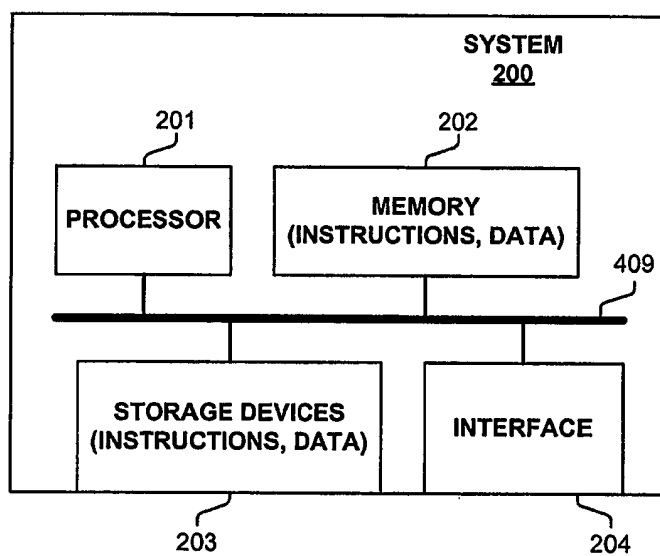


FIGURE 2

6 / 16

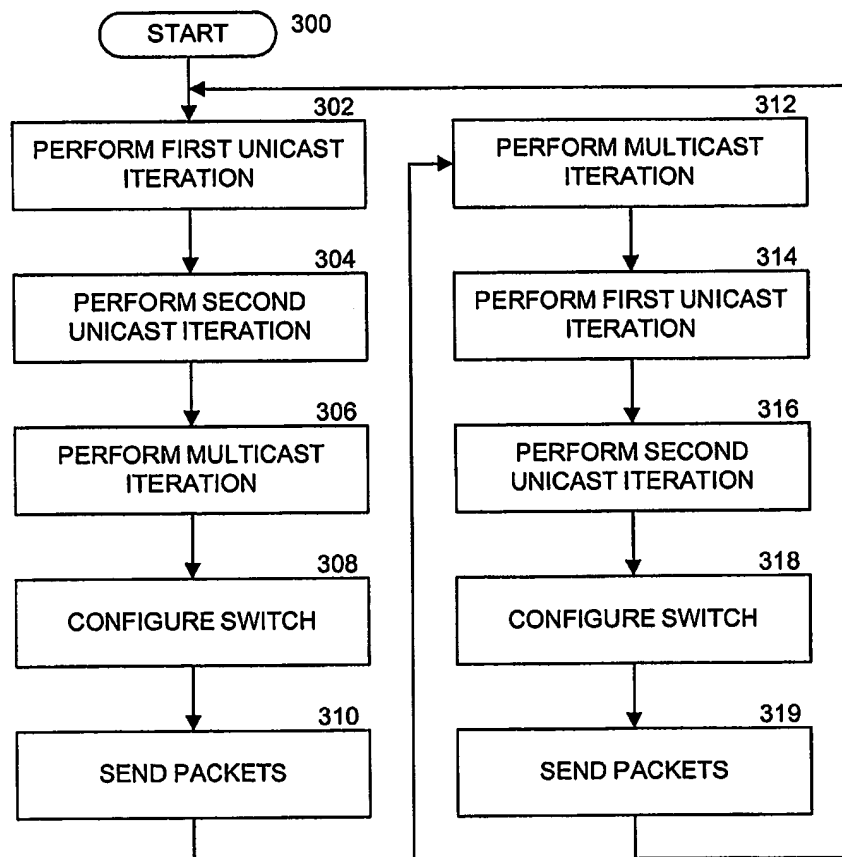
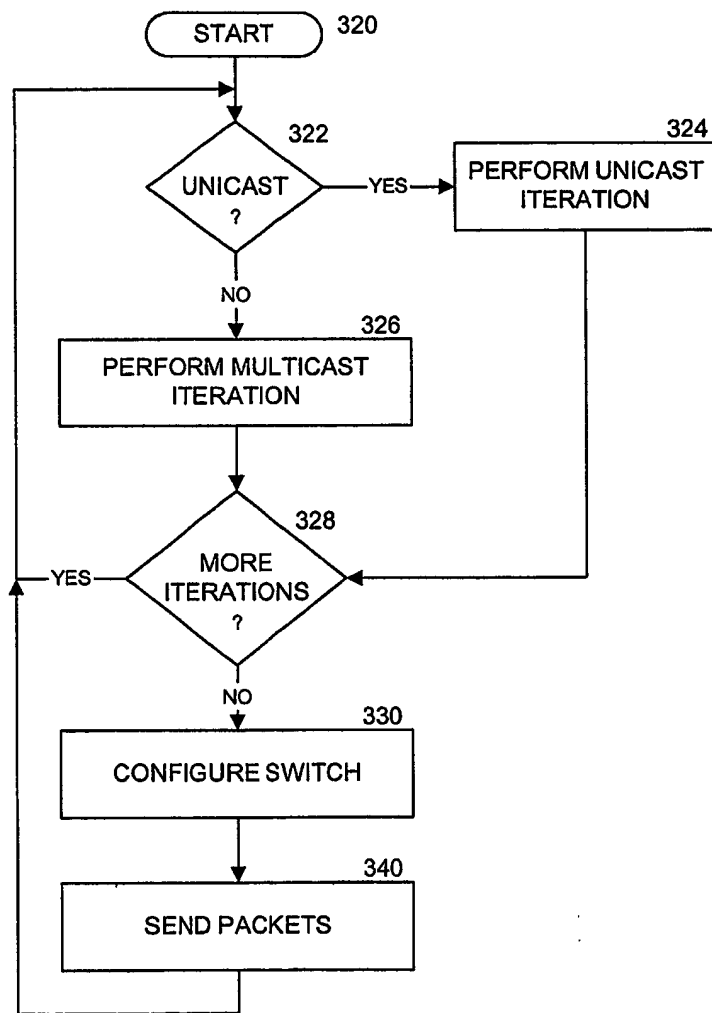
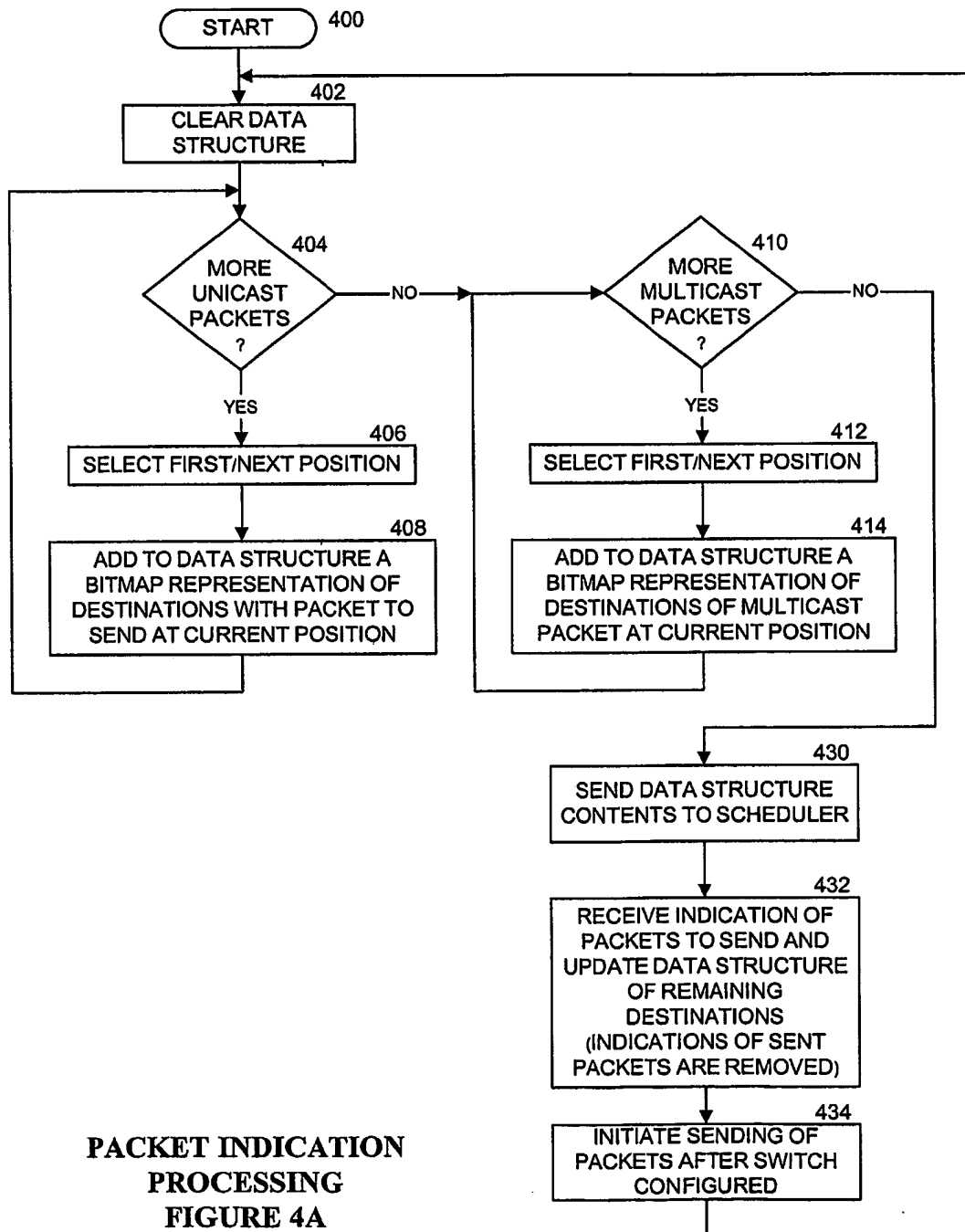


FIGURE 3A

7 / 16

**FIGURE 3B**

8 / 16



9 / 16

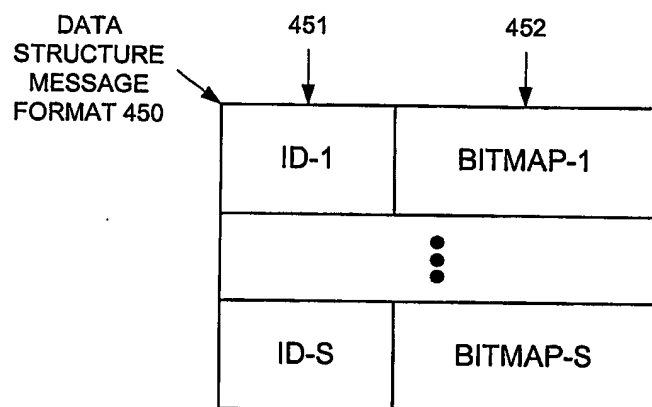


FIGURE 4B

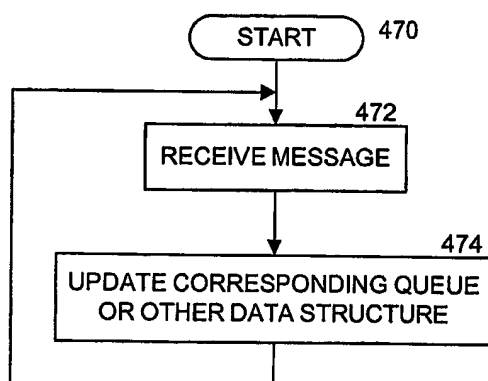
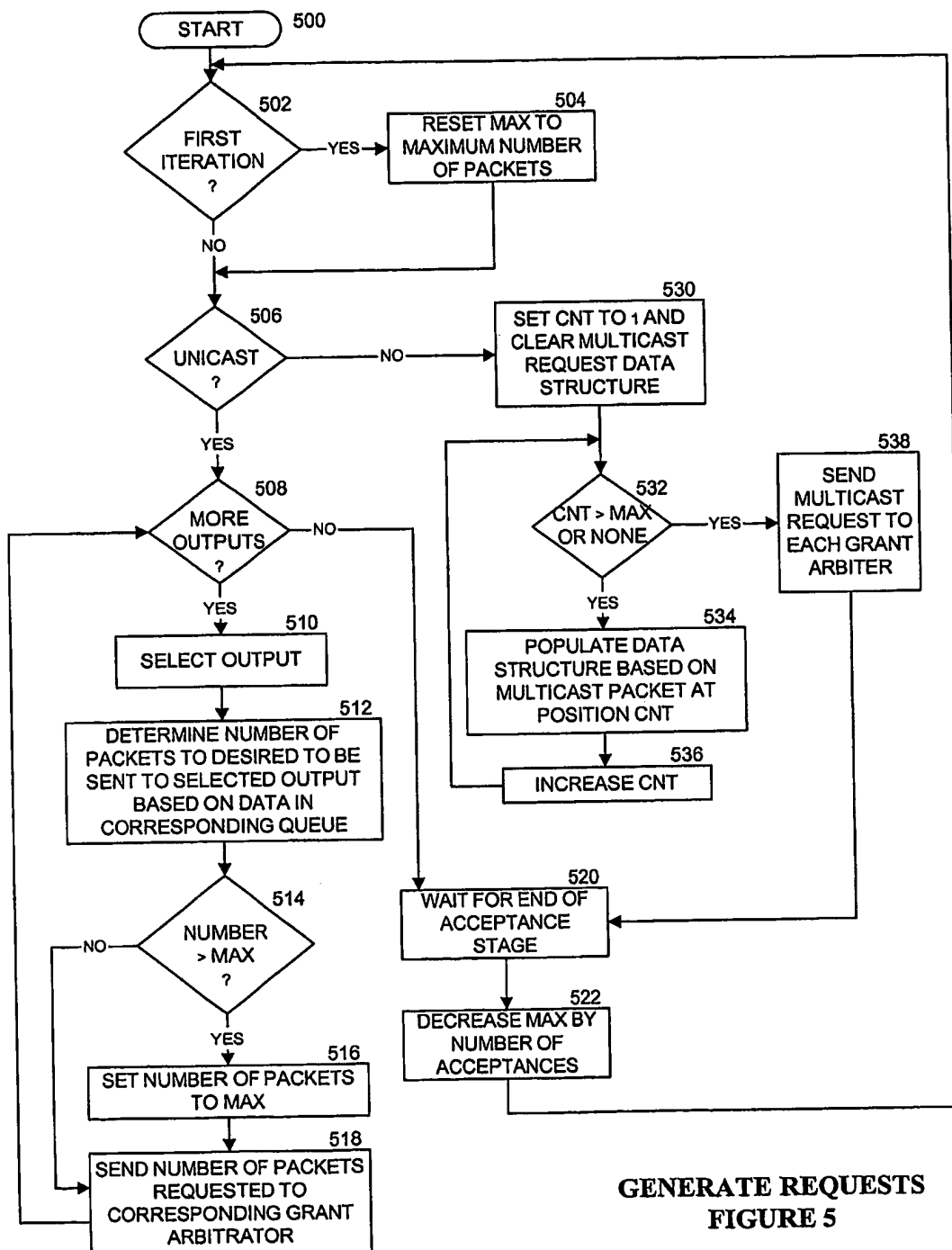
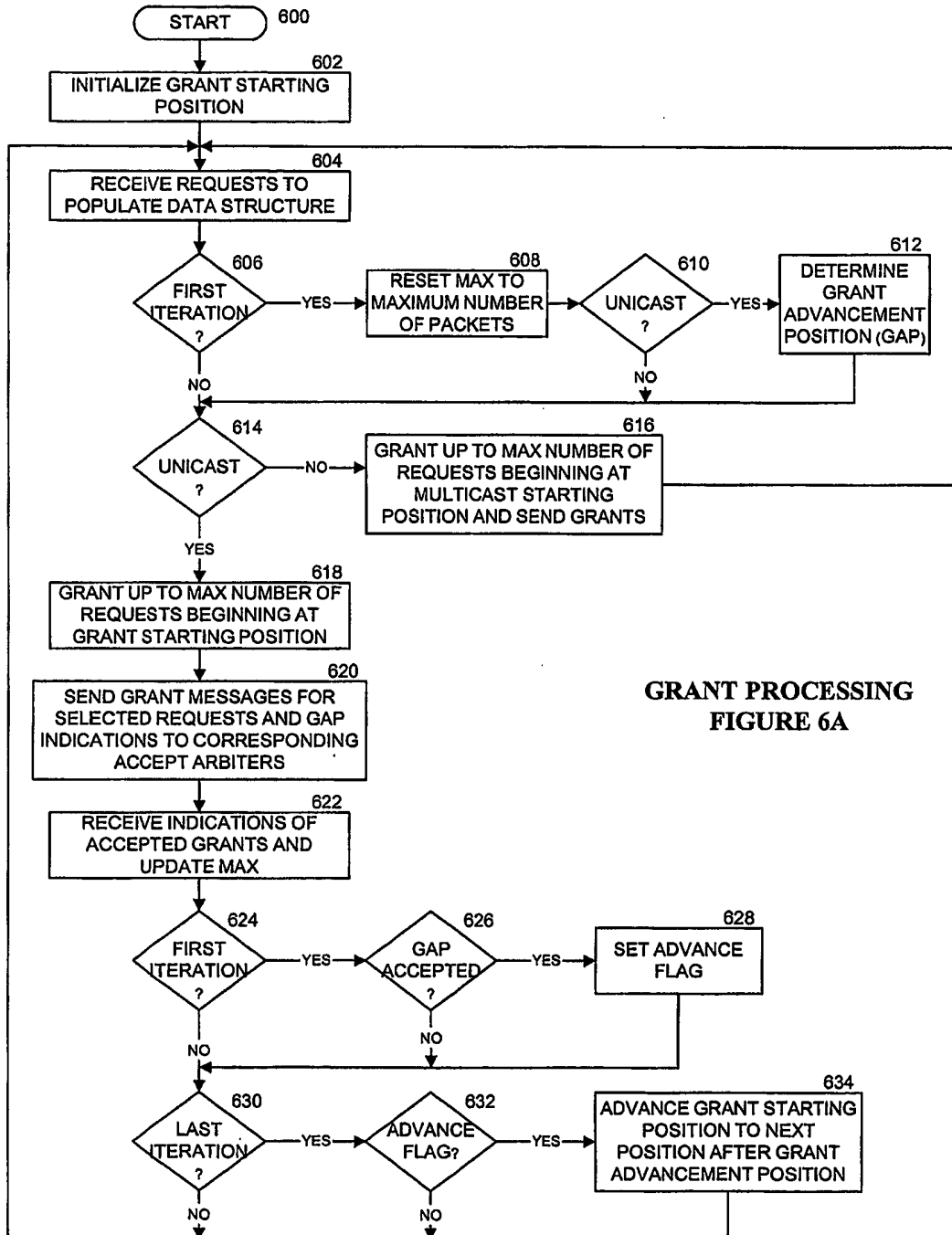


FIGURE 4C

10/16



11 / 16



GRANT PROCESSING
FIGURE 6A

12 / 16

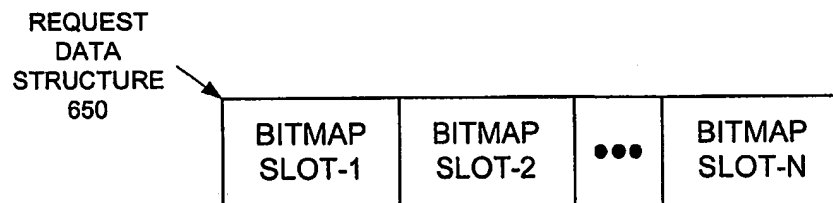


FIGURE 6B

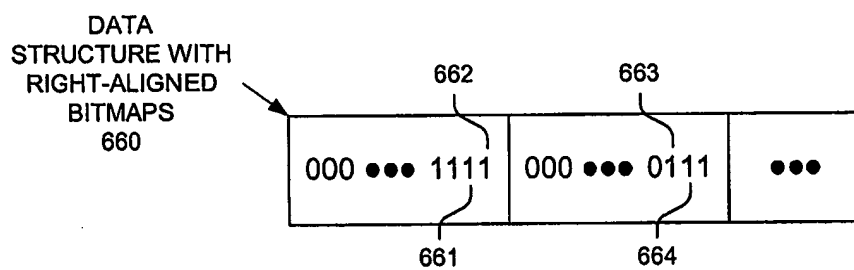
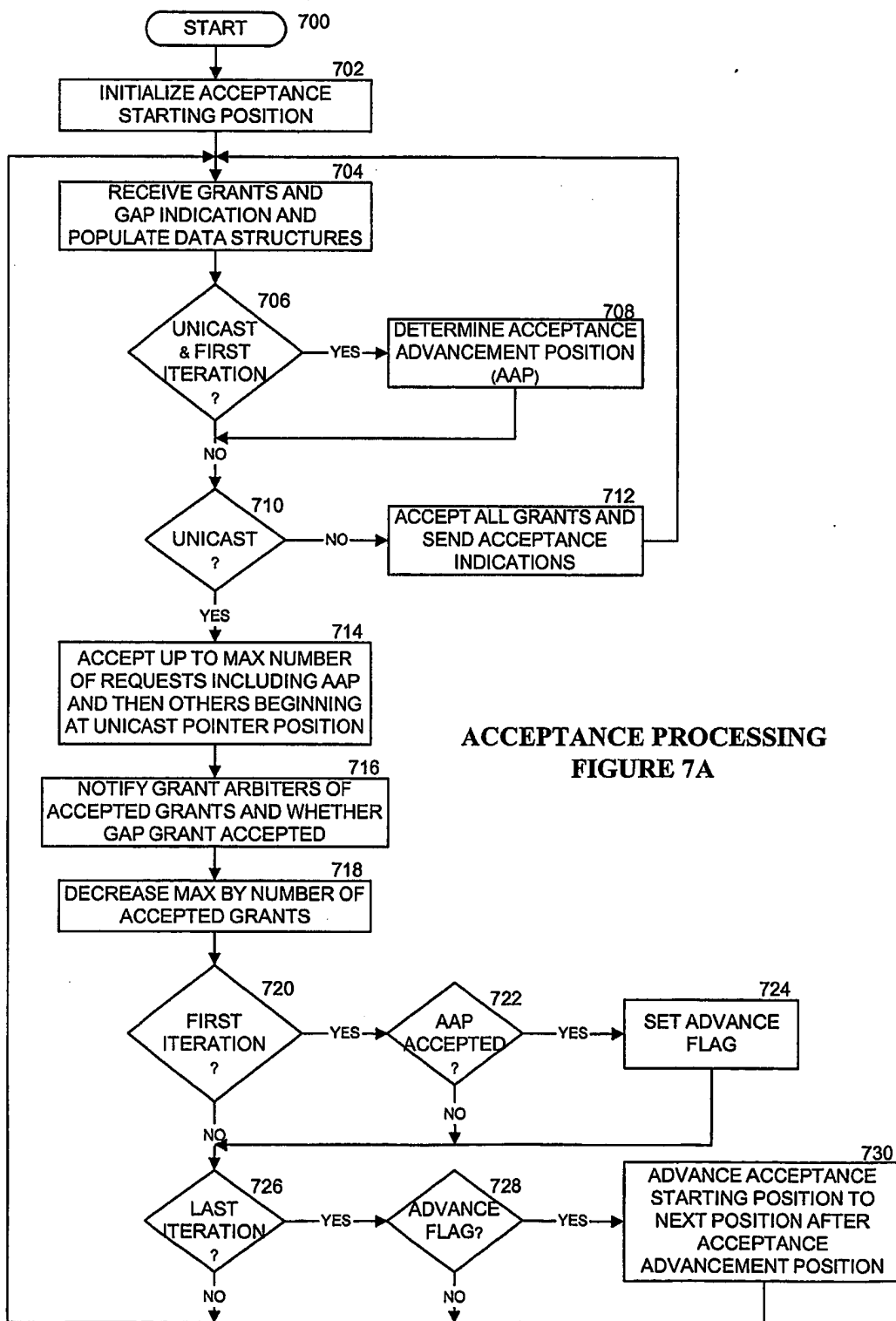
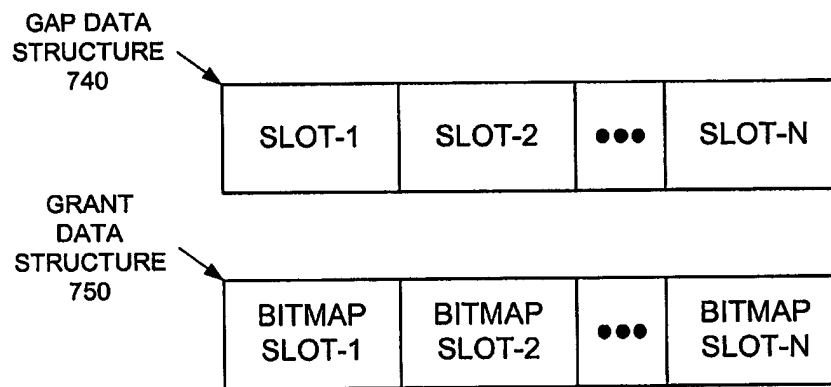
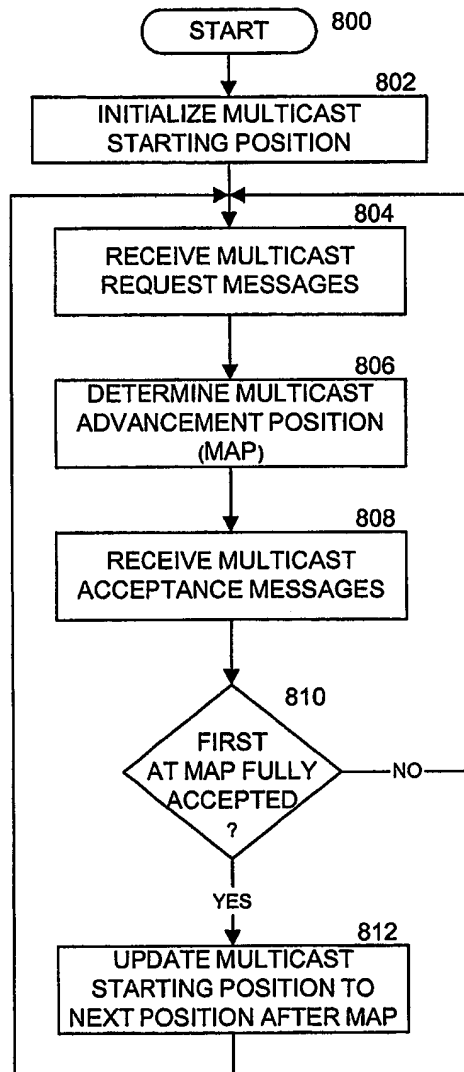


FIGURE 6C



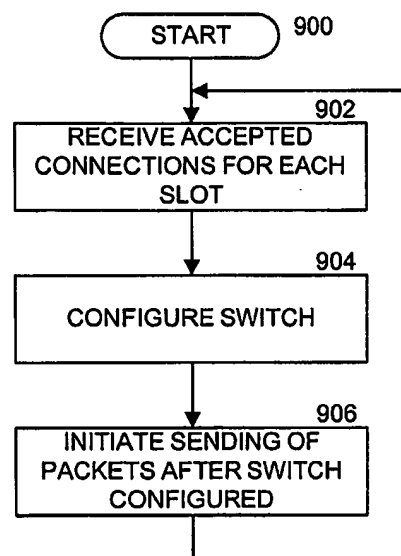
**FIGURE 7B**

15 / 16



**MULTICAST POINTER
PROCESSING
FIGURE 8**

16 / 16

**FIGURE 9**

INTERNATIONAL SEARCH REPORT

Internat Application No
PCT/US 03/08218

A. CLASSIFICATION OF SUBJECT MATTER
IPC 7 H04L12/56 H04Q11/04

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

IPC 7 H04L H04Q

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the International search (name of data base and, where practical, search terms used)

EPO-Internal, WPI Data, INSPEC

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	EP 1 052 814 A (NORTEL NETWORKS LTD) 15 November 2000 (2000-11-15) abstract column 1, line 28-32 column 6, line 32-38 column 9, line 49 -column 13, line 18 column 14, line 28 -column 20, line 35; figures 1,,4,,5,,7A,,7B,,8	1-33
A	WENZHE CUI ET AL: "A threshold based scheduling algorithm for input queue switch" IEEE, 2001, pages 207-212, XP010534273 abstract page 208, left-hand column, line 3 -page 209, left-hand column, line 34; figures 2,3	1-33

☒ Further documents are listed in the continuation of box C.

☒ Patent family members are listed in annex.

* Special categories of cited documents:

A document defining the general state of the art which is not considered to be of particular relevance

E earlier document but published on or after the International filing date

L document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)

O document referring to an oral disclosure, use, exhibition or other means

P document published prior to the International filing date but later than the priority date claimed

T later document published after the International filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention

X document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone

Y document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art.

& document member of the same patent family

Date of the actual completion of the International search

20 August 2003

Date of mailing of the International search report

01/09/2003

Name and mailing address of the ISA

European Patent Office, P.B. 5818 Patentlaan 2
NL - 2280 HV Rijswijk
Tel. (+31-70) 340-2040, Tx. 31 651 epo nl,
Fax: (+31-70) 340-3016

Authorized officer

Fantacone, V

INTERNATIONAL SEARCH REPORT

Intern: Application No
PCT/US 03/08218

C.(Continuation) DOCUMENTS CONSIDERED TO BE RELEVANT		
Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	CHAO H J ET AL: "Centralized contention resolution schemes for a large-capacity optical ATM switch" ATM WORKSHOP PROCEEDINGS, 1998 IEEE FAIRFAX, VA, USA 26-29 MAY 1998, NEW YORK, NY, USA, IEEE, US, 26 May 1998 (1998-05-26), pages 11-16, XP010280499 ISBN: 0-7803-4874-5 the whole document -----	1-33
A	US 2001/043606 A1 (HAN IN-TACK ET AL) 22 November 2001 (2001-11-22) the whole document -----	1-33

INTERNATIONAL SEARCH REPORT

Information on patent family members

Intern: Application No
PCT/US 03/08218

Patent document cited in search report		Publication date	Patent family member(s)	Publication date
EP 1052814	A	15-11-2000	CA 2308243 A1	14-11-2000
			EP 1052814 A2	15-11-2000
			EP 1052815 A2	15-11-2000
<hr/>				
US 2001043606	A1	22-11-2001	NONE	
<hr/>				

**This Page is Inserted by IFW Indexing and Scanning
Operations and is not part of the Official Record**

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images include but are not limited to the items checked:

- ☐ **BLACK BORDERS**
- ☐ **IMAGE CUT OFF AT TOP, BOTTOM OR SIDES**
- ☐ **FADED TEXT OR DRAWING**
- ☐ **BLURRED OR ILLEGIBLE TEXT OR DRAWING**
- ☐ **SKEWED/SLANTED IMAGES**
- ☐ **COLOR OR BLACK AND WHITE PHOTOGRAPHS**
- ☐ **GRAY SCALE DOCUMENTS**
- ☐ **LINES OR MARKS ON ORIGINAL DOCUMENT**
- ☐ **REFERENCE(S) OR EXHIBIT(S) SUBMITTED ARE POOR QUALITY**
- ☐ **OTHER:** _____

IMAGES ARE BEST AVAILABLE COPY.

As rescanning these documents will not correct the image problems checked, please do not report these problems to the IFW Image Problem Mailbox.

THIS PAGE BLANK (USPTO)